



**University of  
Zurich**<sup>UZH</sup>

**Zurich Open Repository and  
Archive**

University of Zurich  
University Library  
Strickhofstrasse 39  
CH-8057 Zurich  
[www.zora.uzh.ch](http://www.zora.uzh.ch)

---

Year: 2018

---

## **Heritability, selection, and the response to selection in the presence of phenotypic measurement error: Effects, cures, and the role of repeated measurements**

Ponzi, Erica ; Keller, Lukas F ; Bonnet, Timothée ; Muff, Stefanie

**Abstract:** Quantitative genetic analyses require extensive measurements of phenotypic traits, a task that is often not trivial, especially in wild populations. On top of instrumental measurement error, some traits may undergo transient (i.e. non-persistent) fluctuations that are biologically irrelevant for selection processes. These two sources of variability, which we denote here as measurement error in a broad sense, are possible causes for bias in the estimation of quantitative genetic parameters. We illustrate how in a continuous trait transient effects with a classical measurement error structure may bias estimates of heritability, selection gradients, and the predicted response to selection. We propose strategies to obtain unbiased estimates with the help of repeated measurements taken at an appropriate temporal scale. However, the fact that in quantitative genetic analyses repeated measurements are also used to isolate permanent environmental instead of transient effects, requires that the information content of repeated measurements is carefully assessed. To this end, we propose to distinguish "short-term" from "long-term" repeats, where the former capture transient variability and the latter the permanent effects. We show how the inclusion of the corresponding variance components in quantitative genetic models yields unbiased estimates of all quantities of interest, and we illustrate the application of the method to data from a Swiss snow vole population.

DOI: <https://doi.org/10.1111/evo.13573>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-157144>

Journal Article

Accepted Version



The following work is licensed under a Creative Commons: Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0) License.

Originally published at:

Ponzi, Erica; Keller, Lukas F; Bonnet, Timothée; Muff, Stefanie (2018). Heritability, selection, and the response to selection in the presence of phenotypic measurement error: Effects, cures, and the role of repeated measurements. *Evolution*, 72(10):1992-2004.

DOI: <https://doi.org/10.1111/evo.13573>

# Heritability, selection, and the response to selection in the presence of phenotypic measurement error: effects, cures, and the role of repeated measurements

Erica Ponzi<sup>1,2</sup>, Lukas F. Keller<sup>1,3</sup>, Timothée Bonnet<sup>1,4</sup>, Stefanie Muff<sup>1,2</sup>

**Running title:** Phenotypic measurement error – effects and cures

<sup>1</sup> Department of Evolutionary Biology and Environmental Studies, University of Zürich, Winterthurerstrasse 190, 8057 Zürich, Switzerland

<sup>2</sup> Department of Biostatistics, Epidemiology, Biostatistics and Prevention Institute, University of Zürich, Hirschengraben 84, 8001 Zürich, Switzerland

<sup>3</sup> Zoological Museum, University of Zürich, Karl-Schmid-Strasse 4, 8006 Zürich, Switzerland

<sup>4</sup> Division of Ecology and Evolution, Research School of Biology, The Australian National University, Acton, Canberra, ACT 2601, Australia

**\* Corresponding author:**

Stefanie Muff,

Department of Evolutionary Biology and Environmental Studies, University of Zürich, Winterthurerstrasse 190, 8057 Zürich, Switzerland

Department of Biostatistics, Epidemiology, Biostatistics and Prevention Institute, University of Zürich, Hirschengraben 84, 8001 Zürich, Switzerland

E-mail: [stefanie.muff@uzh.ch](mailto:stefanie.muff@uzh.ch)

Quantitative genetic analyses require extensive measurements of phenotypic traits, a task that is often not trivial, especially in wild populations. On top of instrumental measurement error, some traits may undergo transient (*i.e.* non-persistent) fluctuations that are biologically irrelevant for selection processes. These two sources of variability, which we denote here as measurement error in a broad sense, are possible causes for bias in the estimation of quantitative genetic parameters. We illustrate how in a continuous trait transient effects with a classical measurement error structure may bias estimates of heritability, selection gradients, and the predicted response to selection. We propose strategies to obtain unbiased estimates with the help of repeated measurements taken at an appropriate temporal scale. However, the fact that in quantitative genetic analyses repeated measurements are also used to isolate permanent environmental instead of transient effects, requires a re-assessment of the information content of repeated measurements. To do so, we propose to distinguish “short-term” from “long-term” repeats, where the former capture transient variability and the latter the permanent effects. We show how the inclusion of the corresponding variance components in quantitative genetic models yields unbiased estimates of all quantities of interest, and we illustrate the application of the method to data from a Swiss snow vole population.

**Keywords:** animal model, breeder’s equation, error variance, permanent environmental effects, quantitative genetics, Robertson-Price identity.

## 52 Introduction

53 Quantitative genetic methods have become increasingly popular for the study of  
 54 natural populations in the last decades, and they now provide powerful tools to in-  
 55 vestigate the inheritance of characters, and to understand and predict evolutionary  
 56 change of phenotypic traits (Falconer and Mackay, 1996; Lynch and Walsh, 1998;  
 57 Charmantier et al., 2014). At its core, quantitative genetics is a statistical approach  
 58 that decomposes the observed phenotype  $P$  into the sum of additive genetic effects  $A$   
 59 and a residual component  $R$ , so that  $P = A + R$ . For simplicity, non-additive genetic  
 60 effects, such as dominance and epistatic effects, are ignored throughout this paper,  
 61 thus the residual component can be thought of as the sum of all environmental ef-  
 62 fects. This basic model can be extended in various ways (Falconer and Mackay, 1996;  
 63 Lynch and Walsh, 1998), with one of the most common being  $P = A + PE + R$ , where  
 64  $PE$  captures *dependent* effects, the so-called *permanent environmental effects*, while  
 65  $R$  captures the residual, *independent* variance that remains unexplained. Permanent  
 66 environmental effects are stable differences among individuals above and beyond the  
 67 permanent differences due to additive genetic effects. In repeated measurements of  
 68 an individual, these effects create within-individual covariation. To prevent inflated  
 69 estimates of additive genetic variance, these effects must therefore be modeled and  
 70 estimated (Lynch and Walsh, 1998; Kruuk, 2004; Wilson et al., 2010).

71 This quantitative genetic decomposition of phenotypes is not possible at the in-  
 72 dividual level in non-clonal organisms, but under the crucial assumption of inde-  
 73 pendence of genetic, permanent environmental, and residual effects, the phenotypic  
 74 variance at the population level can be decomposed into the respective variance  
 75 components as  $\sigma_P^2 = \sigma_A^2 + \sigma_{PE}^2 + \sigma_R^2$ . These variance components can then be used  
 76 to understand and predict evolutionary change of phenotypic traits. For example,  
 77 the additive genetic variance ( $\sigma_A^2$ ) can be used to predict the response to selection  
 78 using the breeder’s equation. It predicts the response to selection  $R_{BE}$  of a trait  $\mathbf{z}$   
 79 (bold face notation denotes vectors) from the product of the heritability ( $h^2$ ) of the  
 80 trait and the strength of selection ( $S$ ) as

$$R_{BE} = h^2 \cdot S \quad (1)$$

81 (Lush, 1937; Falconer and Mackay, 1996), where  $h^2$  is the proportion of additive  
 82 genetic to total phenotypic variance

$$h^2 = \frac{\sigma_A^2}{\sigma_P^2} \quad , \quad (2)$$

83 and  $S$  is the selection differential, defined as the mean phenotypic difference between  
 84 selected individuals and the population mean or, equivalently, the phenotypic covari-

ance  $\sigma_p(\mathbf{z}, \mathbf{w})$  between the trait ( $\mathbf{z}$ ) and relative fitness ( $\mathbf{w}$ ). Besides the breeder's equation, evolution can be predicted using the secondary theorem of selection, according to which evolutionary change is equal to the additive genetic covariance of a trait with relative fitness, that is,

$$R_{\text{STS}} = \sigma_a(\mathbf{z}, \mathbf{w}) \quad (3)$$

(Robertson, 1966; Price, 1970). Morrissey et al. (2010) and Morrissey et al. (2012) discuss the differences between the breeder's equation and the secondary theorem of selection in detail. A major difference is that in contrast to  $R_{\text{BE}}$ ,  $R_{\text{STS}}$  only estimates the population evolutionary trajectory, but does not measure the role of selection in shaping this evolutionary change.

One measure of the role of selection is the selection gradient, which quantifies the strength of natural selection on a trait. For a normally distributed trait ( $\mathbf{z}$ ), it is given as the slope  $\beta_z$  of the linear regression of relative fitness on a phenotypic trait (Lande and Arnold, 1983), that is,

$$\beta_z = \frac{\sigma_p(\mathbf{z}, \mathbf{w})}{\sigma_p^2(\mathbf{z})}, \quad (4)$$

where  $\sigma_p^2(\mathbf{z})$  denotes the phenotypic variance of the trait, for which we only write  $\sigma_p^2$  when there is no ambiguity about what trait the phenotypic variance refers to.

The reliable estimation of the parameters of interest ( $h^2$ ,  $\sigma_p(\mathbf{z}, \mathbf{w})$ ,  $\sigma_a(\mathbf{z}, \mathbf{w})$  and  $\beta_z$ ) and the successful prediction of evolution as  $R_{\text{BE}}$  or  $R_{\text{STS}}$ , require large amounts of data, often collected across multiple generations and with known relationships among individuals in the data set. For many phenotypic traits of interest, data collection is often not trivial, and multiple sources of error, such as phenotypic measurement error, pedigree errors (wrong relationships among individuals), or non-randomly missing data may affect the parameter estimates. Several studies have discussed and addressed pedigree errors (*e.g.* Keller et al., 2001; Griffith et al., 2002; Senneke et al., 2004; Charmantier and Reale, 2005; Hadfield, 2008) and problems arising from missing data (*e.g.* Steinsland et al., 2014; Wolak and Reid, 2017). In contrast, although known for a long time (*e.g.* Price and Boag, 1987), the effects of phenotypic measurement error on estimates of (co-)variance components have received less attention (but see *e.g.* Hoffmann, 2000; Dohm, 2002; Macgregor et al., 2006; van der Sluis et al., 2010; Ge et al., 2017). In particular, general solutions to obtaining unbiased estimates of (co-)variance parameters in the presence of phenotypic measurement error are lacking.

In the simplest case, and the case considered here, phenotypic measurement error is assumed to be independent and additive, that is, instead of the actual phenotype

118  $\mathbf{z}$ , an error-prone version

$$\mathbf{z}^* = \mathbf{z} + \mathbf{e}, \quad \mathbf{e} \sim \mathbf{N}(\mathbf{0}, \sigma_{em}^2 \mathbf{I}) \quad (5)$$

119 is measured, where  $\mathbf{e}$  denotes an error term with independent correlation structure  
 120  $\mathbf{I}$  and error variance  $\sigma_{em}^2$  (see p.121 Lynch and Walsh, 1998). As a consequence,  
 121 the *observed* phenotypic variance of the measured values is  $\sigma_p^2(\mathbf{z}^*) = \sigma_p^2(\mathbf{z}) + \sigma_{em}^2$ ,  
 122 and thus larger than the *actual* phenotypic variance. The error variance  $\sigma_{em}^2$  thus  
 123 must be disentangled from  $\sigma_p^2(\mathbf{z})$  to obtain unbiased estimates of quantitative ge-  
 124 netic parameters. However, most existing methods for continuous trait analyses that  
 125 acknowledged measurement error have modeled it as part of the residual component,  
 126 and thus implicitly as part of the total phenotypic value (*e.g.* Dohm, 2002; Macgre-  
 127 gor et al., 2006; van der Sluis et al., 2010). This means that in the decomposition  
 128 of a phenotype  $P = A + PE + R$ , measurement error is absorbed in  $R$ , thus  $\sigma_{em}^2$   
 129 is absorbed by  $\sigma_R^2$ . This practice effectively *downwardly* biases measures that are  
 130 proportions of the phenotypic variance, in particular  $h^2$  and  $\beta_z$ . To see why, let us  
 131 denote the biased measures as  $h_\star^2$  and  $\beta_z^\star$ . The biased version of heritability is then  
 132 given as

$$h_\star^2 = \frac{\sigma_A^2}{\sigma_P^2 + \sigma_{em}^2} \leq \frac{\sigma_A^2}{\sigma_P^2}, \quad (6)$$

133 because under the assumption taken here that measurement error is independent  
 134 of the actual trait value, measurement error is also independent of additive genetic  
 135 differences and therefore leaves the estimate of the additive genetic variance  $\sigma_A^2$   
 136 unaffected. This was already pointed out *e.g.* by Lynch and Walsh (p.121, 1998) or  
 137 Ge et al. (2017). Equation (6) directly illustrates that  $h_\star^2$  is attenuated by a factor  
 138  $\lambda = \sigma_P^2 / (\sigma_P^2 + \sigma_{em}^2)$ , denoted as reliability ratio (*e.g.* Carroll et al., 2006). Using the  
 139 same argument, one can show that  $\beta_z^\star = \lambda \beta_z$ , but also  $R_{BE}^\star = \lambda R_{BE}$ , as will become  
 140 clear later.

141 To obtain unbiased estimates of  $h^2$ ,  $\beta_z$ , or any other quantity that depends on  
 142 unbiased estimates of  $\sigma_P^2$ , it is thus necessary to disentangle  $\sigma_{em}^2$  from the actual phe-  
 143 notypic variance  $\sigma_P^2$ , and particularly from its residual component  $\sigma_R^2$ . Importantly,  
 144 however, purely mechanistic measurement imprecision is often not the only source  
 145 of variation that may be considered irrelevant for the mechanisms of inheritance and  
 146 selection in the system under study. Here, we therefore follow Ge et al. (2017) and  
 147 use the term “transient effects” for the sum of measurement errors *plus* any biological  
 148 short-term changes of the phenotype itself that are not considered relevant for the  
 149 selection process, briefly denoted as “irrelevant fluctuations” of the actual trait.

150 As an example, if the trait is the mass of an adult animal, repeated measurements  
 151 within the same day are expected to differ even in the absence of instrumental error,

152 simply because animals eat, drink and defecate (for an example of the magnitude  
153 of these effects see Keller and Van Noordwijk, 1993). Such short-term fluctuations  
154 might not be of interest for the study of evolutionary dynamics, if the fluctuations do  
155 not contribute to the selection process in a given population. Under the assumption  
156 that these fluctuations are additive and independent among each other and of the  
157 actual trait value, they are mathematically indistinguishable from pure measurement  
158 error. In the remainder of the paper, we therefore do not introduce a separate  
159 notation to discriminate between (mechanistic) measurement error and biological  
160 short-term fluctuations, but treat them as a single component ( $e$ ) with a total “error”  
161 variance  $\sigma_{e_m}^2$ . Consequently, we may sometimes refer to “measurement error” when  
162 in fact we mean transient effects as the sum of measurement error and transient  
163 fluctuations.

164 The aim of this article is to develop general methods to obtain unbiased estimates  
165 of heritability, selection, and response to selection in the presence of measurement  
166 error and irrelevant fluctuations of a trait, building on the work by Ge et al. (2017).  
167 We start by clarifying the meaning and information content of repeated phenotypic  
168 measurements on the same individual. The type of phenotypic trait we have in  
169 mind is a relatively plastic trait, such as milk production or an animal’s mass, which  
170 are expected to undergo changes across an individual’s lifespan that are relevant  
171 for selection. We show that repeated measures taken over different time intervals  
172 can help separate transient effects from more stable (permanent) environmental and  
173 genetic effects. We proceed to show that based on such a variance decomposition  
174 one can construct models that yield unbiased estimates of heritability, selection, and  
175 the response to selection. We illustrate these approaches with empirical quantitative  
176 trait analyses of body mass measurements taken in a population of snow voles in  
177 the Swiss alps (Bonnet et al., 2017).

## 178 Material and methods

### 179 Short-term and long-term repeated measurements

180 Table 1 gives an overview of how the different parameters considered here are (or  
181 are not) affected by the presence of measurement error. In order to retrieve unbi-  
182 ased estimates of all quantities given in Table 1, we must be able to appropriately  
183 model and estimate the measurement error variance  $\sigma_{e_m}^2$ , which can be achieved  
184 with repeated measurements. These repeated measurements must be taken in close  
185 temporal vicinity, that is, on a time scale where the focal trait is not actually un-  
186 dergoing any phenotypic changes that are relevant for selection. We introduce the  
187 notion of a *measurement session* for such *short-term* time intervals. In other words,



a measurement session can be defined as a sufficiently short period of time during which the investigator is willing to assume that the residual component is constant. On the other hand, measurements are often repeated across much longer periods of time, such as months, seasons, or years, during which phenotypic change is not expected to be solely due to transient effects, and the resulting trait variation is often relevant for selection. Thus, *long-term* repeats, taken across different measurement sessions, help separating permanent environmental effects from residual components (e.g. Wilson et al., 2010).

The distinction between short-term and long-term repeats, and thus the definition of a measurement session, may not always be obvious or unique for a given trait. In the introduction we employed the example of an animal’s mass that transiently fluctuates within a day. Depending on the context, such fluctuations might not be of interest, and the “actual” phenotypic value would correspond to the average daily mass. A reasonable measurement session could then be a single day, and within-day repeats can thus be used to estimate  $\sigma_{e_m}^2$ . If however *any* fluctuations in body mass are of interest, irrespective of how persistent they are, much shorter measurement sessions, such as seconds or minutes, would be appropriate to ensure that only the purely mechanistic measurement error variance is represented by  $\sigma_{e_m}^2$ .

## Repeated measurements in the animal model

In the following we show how measurement error can be incorporated in the key tool of quantitative genetics, the *animal model*, a special type of (generalized) linear mixed model, which is commonly used to decompose the phenotypic variance of a trait into genetic and non-genetic components (Henderson, 1976; Lynch and Walsh, 1998; Kruuk, 2004).

Let us assume that phenotypic measurements of a trait are blurred by measurement error following model (5), and that measurements have been taken both across and within multiple measurement sessions, as indicated in Figure 1a. Denoting by  $z_{ijk}^*$  the  $k^{\text{th}}$  measurement of individual  $i$  in session  $j$ , it is possible to fit a model that decomposes the trait value as

$$z_{ijk}^* = \mu + \mathbf{x}_{ijk}^\top \boldsymbol{\beta} + a_i + id_i + R_{ij} + e_{ijk} , \quad (7)$$

where  $\mu$  is the population intercept,  $\boldsymbol{\beta}$  is a vector of fixed effects and  $\mathbf{x}_{ijk}$  is the vector of covariates for measurement  $k$  in session  $j$  of animal  $i$ . The remaining components are the random effects, namely the breeding value  $a_i$  with dependency structure  $(a_1, \dots, a_n)^T \sim \mathbf{N}(\mathbf{0}, \sigma_A^2 \mathbf{A})$ , an independent, animal-specific permanent environmental effect  $id_i \sim \mathbf{N}(0, \sigma_{PE}^2)$ , an independent Gaussian residual term  $R_{ij} \sim \mathbf{N}(0, \sigma_R^2)$ , and an independent error term  $e_{ijk} \sim \mathbf{N}(0, \sigma_{e_m}^2)$  that absorbs any transient effects



captured by the within-session repeats. The dependency structure of the breeding values  $a_i$  is encoded by the additive genetic relatedness matrix  $\mathbf{A}$  (Lynch and Walsh, 1998), which is traditionally derived from a pedigree, but can alternatively be calculated from genomic data (Meuwissen et al., 2001; Hill, 2014). The model can be further expanded to include more fixed or random effects, such as maternal, nest or time effects, but we omit such terms here without loss of generality. Importantly, model (7) does not require that all individuals have repeated measurements in each session in order to obtain an unbiased estimate of the variance components in the presence of measurement error. In fact, even if there are, on average, fewer than two repeated measurements per individual within sessions, it may be possible to separate the error variance from the residual variance, as long as the total number of within-session repeats over all individuals is reasonably large. We will in the following refer to model (7) as the “error-aware” model.

If, however, a trait has not been measured across different time scales (*i.e.* either only within or only across measurement sessions), not all variance components are estimable. In the first case, when repeats are only taken within a single measurement session for each individual, as depicted in Figure 1b, an error term can be included in the model, but a permanent environmental effect cannot. The model must then be reduced to

$$z_{ik}^* = \mu + \mathbf{x}_{ik}^\top \boldsymbol{\beta} + a_i + R_i + e_{ik} , \quad (8)$$

thus it is possible to estimate the error variance  $\sigma_{e_m}^2$  and to obtain unbiased estimates of  $\sigma_A^2$  and  $h^2$ , while the residual variance  $\sigma_R^2$  then also contains the permanent environmental variance. In the second case, when repeated measurements are only available from across different measurement sessions, as illustrated in Figure 1c, the error variance cannot be estimated. Instead, an animal-specific permanent environmental effect can be added to the model, which is then given as

$$z_{ij}^* = \mu + \mathbf{x}_{ij}^\top \boldsymbol{\beta} + a_i + id_i + R_{ij} \quad (9)$$

for the measurement in session  $j$  for individual  $i$ . Interestingly, this last model mirrors the types of repeats that motivated quantitative geneticists to isolate  $\sigma_{PE}^2$ , which may otherwise be confounded not only with  $\sigma_R^2$ , but also with  $\sigma_A^2$ . This occurs because the repeated measurements across sessions induce an increased within-animal correlation (*i.e.* a similarity) that may be absorbed by  $\sigma_A^2$  if not modeled appropriately (Kruuk and Hadfield, 2007; Wilson et al., 2010).

## Measurement error and selection

Selection occurs when a trait is correlated with fitness, such that variations in the trait values lead to predictable variations among the same individuals in fitness. The leading approach for measuring the strength of directional selection is the one developed by Lande and Arnold (1983), who proposed to estimate the selection gradient  $\beta_z$  as the slope of the regression of relative fitness  $\mathbf{w}$  on the phenotypic trait  $\mathbf{z}$

$$\mathbf{w} = \alpha + \beta_z \cdot \mathbf{z} + \boldsymbol{\epsilon} , \quad (10)$$

with intercept  $\alpha$  and residual error vector  $\boldsymbol{\epsilon}$ . This model can be further extended to account for covariates, such as sex or age. If the phenotype  $\mathbf{z}$  is measured with error (which may again encompass any irrelevant fluctuations), such that the observed value is  $\mathbf{z}^* = \mathbf{z} + \mathbf{e}$  with error variance  $\sigma_{e_m}^2$  as in (5), the regression of  $\mathbf{w}$  against  $\mathbf{z}^*$  leads to an attenuated version of  $\beta_z$  (Mitchell-Olds and Shaw, 1987; Fuller, 1987; Carroll et al., 2006). Using that  $\hat{\beta}_z = \frac{\sigma_p(\mathbf{z}, \mathbf{w})}{\sigma_p^2(\mathbf{z})}$ ,  $\sigma_p^2(\mathbf{z}^*) = \sigma_p^2(\mathbf{z}) + \sigma_{e_m}^2$ , and the assumption that the error in  $\mathbf{z}^*$  is independent of  $\mathbf{w}$ , simple calculations show that the error-prone estimate of selection is

$$\hat{\beta}_z^* = \frac{\sigma_p(\mathbf{z}^*, \mathbf{w})}{\sigma_p^2(\mathbf{z}^*)} = \frac{\sigma_p(\mathbf{z}, \mathbf{w})}{\sigma_p^2(\mathbf{z}) + \sigma_{e_m}^2} \leq \hat{\beta}_z .$$

Hence, the quantity that is estimated is  $\beta_z^* = \lambda \beta_z$  with  $\lambda = \sigma_p^2(\mathbf{z}) / (\sigma_p^2(\mathbf{z}) + \sigma_{e_m}^2)$ , thus  $\beta_z$  suffers from exactly the same bias as the estimate of heritability (see again Table 1). To obtain an unbiased estimate of selection it may thus often be necessary to account for the error by a suitable error model. Such error-aware model must rely on the same type of short-term repeated measurements as those used in (7) or (8), but with the additional complication that  $\mathbf{z}$  is now a covariate in a regression model, and no longer the response. In order to estimate an unbiased version of  $\beta_z$  we therefore rely on the interpretation as an error-in-variables problem for classical measurement error (Fuller, 1987; Carroll et al., 2006). To this end, we propose to formulate a *Bayesian hierarchical model*, because this formulation, together with the possibility to include prior knowledge, provides a flexible way to model measurement error (Stephens and Dellaportas, 1992; Richardson and Gilks, 1993). To obtain an error-aware model that accounts for error in selection gradients, we need a three-level hierarchical model: The first level is the regression model for selection, and the second level is given by the error model of the observed covariate  $\mathbf{z}^*$  given its true value  $\mathbf{z}$ . Third, a so-called *exposure model* for the unobserved (true) trait value is required to inform the model about the distribution of  $\mathbf{z}$ , and it seems natural to employ the animal model (9) for this purpose. Again using the notation for an individual  $i$  measured in different sessions  $j$  and with repeats  $k$  within sessions, the

288 formulation of the three-level hierarchical model is given as

$$w_{ij} = \alpha + \beta_z z_{ij} + \mathbf{x}_{ij}^\top \boldsymbol{\beta} + \epsilon_{ij}, \quad \epsilon_{ij} \sim \mathbf{N}(0, \sigma_\epsilon^2) \quad \text{Selection model} \quad (11)$$

$$z_{ijk}^* = z_{ij} + e_{ijk}, \quad e_{ijk} \sim \mathbf{N}(0, \sigma_{e_m}^2) \quad \text{Error model} \quad (12)$$

$$z_{ij} = \mu + \mathbf{x}_{ij}^\top \boldsymbol{\beta} + a_i + id_i + R_{ij}, \quad R_{ij} \sim \mathbf{N}(0, \sigma_R^2) \quad \text{Exposure model} \quad (13)$$

289 where  $w_{ij}$  is the measurement of relative fitness for individual  $i$ , usually taken only  
 290 once per individual and having the same value for all measurement sessions  $j$ ,  $\boldsymbol{\beta}$  is  
 291 a vector of fixed effects,  $\mathbf{x}_{ij}$  is the vector of covariates for animal  $i$  in measurement  
 292 session  $j$ ,  $\beta_z$  is the selection gradient, and  $\alpha$  and  $\epsilon_{ij}$  are respectively the intercept  
 293 and the independent residual term from the linear regression model. The classical  
 294 independent measurement error term is given by  $e_{ijk}$ . This formulation as a hier-  
 295 archical model gives an unbiased estimate of the selection gradient  $\beta_z$ , because the  
 296 lower levels of the model properly account for the error in  $\mathbf{z}$  by explicitly modelling  
 297 it. It might be helpful to see that the second and third levels are just a hierarchical  
 298 representation of model (7). Model (11)-(13) can be fitted in a Bayesian setup, see  
 299 for instance Muff et al. (2015) for a description of the implementation in INLA (Rue  
 300 et al., 2009) via its R interface R-INLA.

301 Note that model (11) is formulated here for directional selection. Although the  
 302 explicit discussion of alternative selection mechanisms, such as stabilizing or disrup-  
 303 tive selection, is beyond the scope of the present paper, we note that error modelling  
 304 for these cases is straightforward: The only change is that the linear selection model  
 305 (11) is replaced by the appropriate alternative, for example by including quadratic  
 306 or any other kind of non-linear terms (*e.g.* Fisher, 1930; Lande and Arnold, 1983).  
 307 Moreover, (11) can be replaced by any other regression model, for example by one  
 308 that accounts for non-normality of fitness (see *e.g.* Morrissey and Sakrejda, 2013;  
 309 Morrissey and Goudie, 2016). Similarly, it is conceptually straightforward to replace  
 310 the Gaussian error and exposure models, if there is reason to believe that the normal  
 311 assumptions for the error term  $e_{ijk}$  or the residual term  $R_{ij}$  are unrealistic, for ex-  
 312 ample if  $\mathbf{z}$  is a count or a binary variable. In fact, equation (10) to estimate selection  
 313 does not actually assume a specific distribution for  $\mathbf{z}$ , however the interpretation of  
 314  $\beta_z$  as a directional selection gradient to predict evolutionary change may be lost for  
 315 non-Gaussian traits (Lande and Arnold, 1983). Finally and importantly, although  
 316 multivariate selection is not covered in the present paper, it is possible to extend  
 317 the hierarchical model (11)-(13) to the multivariate case.

## Measurement error and the response to selection

### The breeder's equation

Evolutionary response to a selection process on a phenotypic trait can be predicted either by the breeder's equation (1) or by the Robertson-Price identity (3), and these two approaches are equivalent only when the respective trait value (in the univariate model) is the sole causal factor affecting fitness (Morrissey et al., 2010, 2012). Even if the breeder's equation is formulated for multiple traits, the implicit assumption still is that *all* correlated traits causally related to fitness are included in the model. Given that fitness is a complex trait that usually depends on many unmeasured variables (Møller and Jennions, 2002; Peek et al., 2003), it is not surprising that the breeder's equation is often not successful in predicting evolutionary change in natural systems (Hadfield, 2008; Morrissey et al., 2010), in contrast to (artificial) animal breeding situations, where, thanks to the control over the process, all the traits affecting fitness are known and included in the models (Lush, 1937; Falconer and Mackay, 1996; Roff, 2007).

To understand how transient effects affect the estimate of  $R_{BE} = h^2 \cdot S$ , we must understand how the components  $h^2$  and  $S$  are affected. We have seen that  $h_\star^2 = \lambda h^2$ . On the other hand, the selection differential  $S^\star = \sigma_p(\mathbf{z}^\star, \mathbf{w})$  is an unbiased estimate of  $\sigma_p(\mathbf{z}, \mathbf{w})$ , because under the assumption of independence of the error vector  $\mathbf{e}$  and fitness  $\mathbf{w}$ ,

$$\sigma_p(\mathbf{z}^\star, \mathbf{w}) = \sigma_p(\mathbf{z} + \mathbf{e}, \mathbf{w}) = \sigma_p(\mathbf{z}, \mathbf{w}) + \underbrace{\sigma_p(\mathbf{e}, \mathbf{w})}_{=0} = \sigma_p(\mathbf{z}, \mathbf{w}) . \quad (14)$$

Consequently, the bias in  $h_\star^2$  directly propagates to the estimated response to selection, that is,  $R_{BE}^\star = \lambda R_{BE}$  (Table 1).

### The Robertson-Price identity

Response to selection can also be predicted using the secondary theorem of selection. Specifically, the additive genetic covariance of the relative fitness  $\mathbf{w}$  and the phenotypic trait  $\mathbf{z}$ ,  $\sigma_a(\mathbf{w}, \mathbf{z})$  can be estimated from a bivariate animal model. If interest centers around the evolutionary response of a single trait, the model for the response vector including the (error-prone) trait values  $\mathbf{z}^\star$  and relative fitness values  $\mathbf{w}$  is bivariate with

$$\begin{bmatrix} \mathbf{z}^\star \\ \mathbf{w} \end{bmatrix} = \boldsymbol{\mu} + \mathbf{X}\boldsymbol{\beta} + \mathbf{D}\mathbf{a} + \mathbf{Z}\mathbf{r} , \quad (15)$$

where  $\boldsymbol{\mu}$  is the intercept vector,  $\boldsymbol{\beta}$  the vector of fixed effects,  $\mathbf{X}$  the corresponding design matrix,  $\mathbf{D}$  is the design matrix for the breeding values  $\mathbf{a}$ , and  $\mathbf{Z}$  is a design

matrix for additional random terms  $\mathbf{r}$ . These may include environmental and/or error terms, depending on the structure of the data, that may correspond to the univariate cases of equations (7) - (9) or again to other random terms such as maternal or nest effects. The actual component of interest is the vector of breeding values, which is assumed multivariate normally distributed with

$$\mathbf{a} = \begin{bmatrix} \mathbf{a}(z^*) \\ \mathbf{a}(w) \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \sigma_a^2(z^*)\mathbf{A} & \sigma_a(w, z^*)\mathbf{A} \\ \sigma_a(w, z^*)\mathbf{A} & \sigma_a^2(w)\mathbf{A} \end{bmatrix} \right), \quad (16)$$

where  $\mathbf{a}(z^*)$  and  $\mathbf{a}(w)$  are the respective subvectors for the trait and fitness, and  $\mathbf{A}$  is the relationship matrix derived from the pedigree. An estimate of the additive genetic covariance  $\sigma_a(w, z^*)$  is extracted from this covariance matrix. An interesting feature of the additive genetic covariance, and consequently estimates of the response to selection using the STS, is that it is unbiased by independent error in the phenotype. This can be seen by reiterating the exact same argument as in equation (14), but replacing the phenotypic with the genetic covariance.

We confirmed all these theoretical expectations with a simulation study, where we analysed the effects of measurement error on the estimates of interest by adding error terms with different variances to the phenotypic traits. Details and results of the simulations are given in Appendix 2, while the code for their implementation is reported in Appendix 3.

## Example: Body mass of snow voles

The empirical data we use here stem from a snow vole population that has been monitored between 2006 and 2014 in the Swiss Alps (Bonnet et al., 2017). The genetic pedigree is available for 937 voles, together with measurements on morphological and life history traits. Thanks to the isolated location, it was possible to monitor the whole population and to obtain high recapture probabilities ( $0.924 \pm 0.012$  for adults and  $0.814 \pm 0.030$  for juveniles). Details of the study are given in Bonnet et al. (2017).

Our analyses focused on the estimation of quantitative genetic parameters for the animals' body mass (in grams). The dataset contained 3266 mass observations from 917 different voles across 9 years. Such measurements are expected to suffer from classical measurement error, as they were taken with a spring scale, which is prone to measurement error under field conditions. In addition, the actual mass of an animal may contain irrelevant within-day fluctuations (eating, defecating, digestive processes), but also unknown pregnancy conditions in females, which cannot reliably be determined in the field. Repeated measurements were available, both recorded within and across different seasons. In each season two to five "trapping sessions"

were conducted, which each lasted four consecutive nights. Although this definition of measurement session was based purely on operational aspects driven by the data collection process, we used this time interval to estimate  $\sigma_{e_m}^2$ . It is arguably possible that four days might be undesirably long, and that variability in such an interval includes more than purely transient effects, but the data did not allow for a finer time-resolution. However, to illustrate the importance of the measurement session length, we also repeated all analyses with measurement sessions defined as a calendar month, which is expected to identify a larger (and probably too high) proportion of variance as  $\sigma_{e_m}^2$ . The number of 4-day measurement sessions per individual was on average 3.02 (min = 1, max = 24) with 1.15 (min = 1, max = 3) number of short-term repeats on average, while there were 2.37 (min = 1, max = 13) one-month measurement sessions on average, with 1.41 (min = 1, max = 6) short-term repeats per measurement session.

## Heritability

Bonnet et al. (2017) estimated heritability using an animal model with sex, age, Julian date (JD), squared Julian date and the two-way and three-way interactions among sex, age and Julian date as fixed effects. The inbreeding coefficient was included to avoid bias in the estimation of additive genetic variances (de Boer and Hoeschele, 1993). The breeding value ( $a_i$ ), the maternal identity ( $m_i$ ) and the permanent environmental effect explained by the individual identity ( $id_i$ ) were included as individual-specific random effects.

If no distinction is made between short-term (within measurement session) and long-term (across measurement sessions) repeated measurements, the model that we denote as the *naïve* model is given as

$$z_{ijk}^* = \mu + \mathbf{x}_{ijk}^\top \boldsymbol{\beta} + a_i + m_i + id_i + R_{ijk}, \quad (17)$$

where  $z_{ijk}^*$  is the mass of animal  $i$  in measurement session  $j$  for repeat  $k$ . This model is prone to underestimate heritability, because it does not separate the variance  $\sigma_{e_m}^2$  from the residual variability, and  $\sigma_{e_m}^2$  is thus treated as part of the total phenotypic trait variability. To isolate the measurement error variance, the model expansion

$$z_{ijk}^* = \mu + \mathbf{x}_{ijk}^\top \boldsymbol{\beta} + a_i + m_i + id_i + R_{ij} + e_{ijk},$$

with  $R_{ij} \sim \mathbf{N}(0, \sigma_R^2)$  and  $e_{ijk} \sim \mathbf{N}(0, \sigma_{e_m}^2)$  leads to what we denote here as the *error-aware* model. Under the assumption that the length of a measurement session was defined in an appropriate way, and that the error obeys model (5), this model yields an unbiased estimate of  $h^2$ , calculated as  $\frac{\sigma_A^2}{\sigma_A^2 + \sigma_M^2 + \sigma_{PE}^2 + \sigma_R^2}$  (in agreement with

415 Bonnet et al., 2017), where  $\sigma_{e_m}^2$  is explicitly estimated and thus not included in  
 416 the denominator. Both models were implemented in `MCMCglmm` and are reported  
 417 in Appendix 4. Inverse gamma priors  $\text{IG}(0.01, 0.01)$ , parameterized with shape and  
 418 rate parameters, were used for all variances in all models, while  $\text{N}(0, 10^{12})$  (*i.e.*  
 419 default `MCMCglmm`) priors were given to the fixed effect parameters. Analyses were  
 420 repeated with varying priors on  $\sigma_{e_m}^2$  for a sensitivity check, but results were very  
 421 robust (results not shown).

## 422 Selection

423 Selection gradients were estimated from the regression of relative fitness ( $w$ ) on body  
 424 mass ( $z^*$ ). Relative fitness was defined as the relative lifetime reproductive success  
 425 (rLRS), calculated as the number of offspring over the lifetime of an individual,  
 426 divided by the population mean LRS. The naive estimate of the selection gradient  
 427 was obtained from a linear mixed model (*i.e.* treating rLRS as continuous trait),  
 428 where body mass, sex and age were included as fixed effects, plus a cohort-specific  
 429 random effect. The error-aware version of the selection gradient  $\beta_z$  was estimated  
 430 using a three-layer hierarchical error model as in (11)-(13) that also included an  
 431 additional random effect for cohort in the regression model. Sex and age were also  
 432 included as fixed effects in the exposure model, plus breeding values, permanent  
 433 environmental and a residual term as random effects. The hierarchical model used  
 434 to estimate the error-aware  $\beta_z$  was implemented in `INLA` and is described in Appendix  
 435 1, with R code given in Appendix 5. Again,  $\text{IG}(0.01, 0.01)$  priors were assigned to  
 436 all variance components, while independent  $\text{N}(0, 10^2)$  priors were used for all slope  
 437 parameters. Since rLRS is not actually a Gaussian trait,  $p$ -values and CIs of the  
 438 estimate for  $\beta_z$  from the linear regression model are, however, incorrect. Although  
 439 recent considerations indicate that selection gradients could directly be extracted  
 440 from an overdispersed Poisson model (Morrissey and Goudie, 2016), we followed  
 441 the original analysis of Bonnet et al. (2017) and extracted  $p$ -values from an over-  
 442 dispersed Poisson regression model with absolute LRS as a count outcome, both  
 443 for the (naive) model without error modelling *and* for the hierarchical error model,  
 444 where the linear model (13) was replaced by an overdispersed Poisson regression  
 445 model (see Appendices 1 and 5 for the model description and code for both models).

## 446 Response to selection

447 Response to selection on body mass was estimated with rLRS using the breeder's  
 448 equation (1) and the secondary theorem of selection (3), both for the naive and  
 449 the error-aware versions of the model. The naive and error-aware versions of  $R_{\text{BE}}$   
 450 were estimated by substituting either the naive  $h_{\star}^2$  or the error-aware estimates of



451  $h^2$  into the breeder's equation, where the selection differential was calculated as  
 452 the phenotypic covariance between mass and rLRS. On the other hand,  $R_{\text{STS}}$  was  
 453 estimated from the bivariate animal model, implemented in `MCMCglmm` using the  
 454 same fixed and random effects as those in equation (17). Again  $\text{IG}(0.01, 0.01)$  priors  
 455 were used for the variance components. No residual component was included for the  
 456 fitness trait, as suggested by Morrissey et al. (2012), and its error variance was fixed  
 457 at 0, because no error modelling is required. Appendix 6 contains the respective R  
 458 code.

## 459 Results

### 460 Heritability

461 As expected from theory (Table 1), transient effects in the measurements of body  
 462 mass biased some, but not all, quantitative genetic estimates in our snow vole exam-  
 463 ple (Table 2). The estimates and confidence intervals of the additive genetic variance  
 464  $\sigma_A^2$ , as well as the permanent environmental variance  $\sigma_{PE}^2$  and the maternal variance  
 465 (denoted as  $\sigma_M^2$ ) were only slightly corrected in the error-aware models. Residual  
 466 variances, however, were much lower when measurement error was accounted for in  
 467 the models. The measurement error model separated residual and transient (error)  
 468 variance so that  $\hat{\sigma}_R^2 + \hat{\sigma}_{em}^2$  corresponded approximately to  $\hat{\sigma}_R^2$  from the naive model.  
 469 The overestimation of the residual variance resulted in estimates of heritability that  
 470 were underestimated by nearly 40% when measurement error was ignored ( $\hat{h}^2 = 0.14$   
 471 in the naive model and  $\hat{h}^2 = 0.23$  in the error-aware model).

472 As expected, the estimated measurement error variance is larger when a mea-  
 473 surement session is defined as a full month ( $\hat{\sigma}_{em}^2 = 7.91$ ) than as a 4-day interval  
 474 ( $\hat{\sigma}_{em}^2 = 6.07$ , Table 2), because the trait then has more time and opportunity to  
 475 change. As a consequence, heritability is even slightly higher ( $\hat{h}^2 = 0.24$ ) when the  
 476 longer measurement session definition is used. This example is instructive because it  
 477 underlines the importance of defining the time scale at which short-term repeats are  
 478 expected to capture only transient, and not biologically relevant variability of the  
 479 phenotypic trait. In the case of the mass of a snow vole, most biologists would prob-  
 480 ably agree that changes in body mass over a one-month measurement session may  
 481 well be biologically meaningful (*i.e.* body fat accumulation, pregnancy in females,  
 482 etc.), while it is less clear how much of the fluctuations within a 4-day measurement  
 483 session are transient, and what part of it would be relevant for selection. Within-  
 484 day repeats might be the most appropriate for the case of mass, since within-day  
 485 variance is likely mostly transient, but because the data were not collected with the  
 486 intention to quantify such effects, within-day repeats were not available in sufficient

487 numbers in our example data set.

## 488 Selection

489 As expected, estimates of selection gradients ( $\hat{\beta}_z$ ) obtained with the measurement  
490 error models provided nearly 40% higher estimates of selection than the naive model  
491 (Table 3). The two measurement session lengths yielded similar results. With  
492 and without measurement error modelling, the  $p$ -values of the zero-inflated Poisson  
493 models confirmed the presence of selection on body mass in snow voles ( $p < 0.001$   
494 in all models).

## 495 Response to selection

496 In line with theory, estimates of the response to selection using the breeder's equation  
497 were nearly 40% higher when transient effects were incorporated in the quantitative  
498 genetic models using 4-day measurement sessions ( $\hat{R}_{BE} = 0.10$  in the naive model  
499 and  $\hat{R}_{BE} = 0.16$  in the error-aware model; Table 4). As in the case of heritability, the  
500 one-month measurement session definition resulted in even slightly higher estimates  
501 of the response to selection ( $\hat{R}_{BE} = 0.17$ ). In contrast, response to selection mea-  
502 sured by the secondary theorem of selection  $\hat{R}_{STS}$  did not show evidence of bias, and  
503 the error-aware model with a 4-day measurement session definition estimated the  
504 same value ( $\hat{R}_{STS} = -0.17$ ) as the naive model (Table 4). With a one-month mea-  
505 surement session, we obtained a slightly attenuated value ( $\hat{R}_{STS} = -0.14$ ), although  
506 the difference was small in comparison to the credible intervals (Table 4).

507 This example illustrates that the breeder's equation is generally prone to under-  
508 estimation of the selection response in real study systems when measurement error  
509 in the phenotype is present (Table 1). The results also confirm that estimates for  
510 response to selection may differ dramatically between the breeder's equation and the  
511 secondary theorem of selection approach. As already noticed by Bonnet et al. (2017),  
512 the predicted evolutionary response derived from the breeder's equation points in  
513 the opposite direction in the snow vole data than the estimate derived from the  
514 secondary theorem of selection (*e.g.* naive estimates  $\hat{R}_{BE} = 0.10$  vs.  $\hat{R}_{STS} = -0.17$ ,  
515 with non-overlapping credible intervals; Table 4).

## 516 Discussion

517 This study addresses the problem of measurement error and transient fluctuations  
518 in continuous phenotypic traits in quantitative genetic analyses. We show that mea-  
519 surement error and transient fluctuations can lead to substantial bias in estimates of  
520 several important quantitative genetic parameters, including heritability, selection

521 gradients and the response to selection (Table 1). We introduce modelling strategies  
522 to obtain unbiased estimates in these parameters in the presence of measurement  
523 error and transient fluctuations. These strategies rely on the distinction between  
524 variability from stable effects that are part of the biologically relevant phenotypic  
525 variability, and transient effects, which are the sum of mechanistic measurement er-  
526 ror and biological fluctuations that are considered irrelevant for the selection process.  
527 We argue that ignoring the distinction between stable and transient effects may not  
528 only lead to an *underestimation* of the heritability due to inflated estimates of the  
529 residual variance,  $\sigma_R^2$ , but also to bias in the estimates of selection gradients and the  
530 response to selection. Measurements of the same individual repeated at appropriate  
531 time scales allow the variance from such transient effects to be partitioned, and thus  
532 prevent such bias.

533 How can repeated measurements be used to prevent an *underestimation* of her-  
534 itability, selection, and response to selection, while permanent environment effects  
535 are required in quantitative genetic models of repeated measures to avoid an upward  
536 bias of  $\sigma_A^2$  and, hence, an *overestimation* of  $h^2$  (Wilson et al., 2010)? The fact that  
537 repeated measurements are used to prevent opposite biases in heritability estimates  
538 makes it apparent that the information content in what is termed “repeated measure-  
539 ments” in both cases is very different. The crucial aspect is that it matters at which  
540 temporal distance the repeats were taken, and that the relevance of this distance  
541 depends on the kind of trait under study. Repeats taken on the same individual  
542 at different life stages (“long-term” repeats, *e.g.* across what we call measurement  
543 sessions here) can be used to separate the animal-specific permanent environmental  
544 effect from both genetic and residual variances. On the other hand, repeats taken  
545 in temporal vicinity (“short-term” repeats, *e.g.* within a measurement session) help  
546 disentangle any transient from the residual effects. Only by modelling *both* types of  
547 repeats, that is, across different relevant time scales, is it practically feasible to sep-  
548 arate all variance components. To do so, the quantitative genetic model for the trait  
549 value, typically the animal model, needs extension to three levels of measurement  
550 hierarchy (equation (7)): the individual ( $i$ ), the measurement session ( $j$  within  $i$ )  
551 and the repeat ( $k$  within  $j$  within  $i$ ). As highlighted with the snow vole example, it  
552 may not always be trivial to determine, in a particular system, an appropriate dis-  
553 tinction between short-term and long-term repeats, and consequently how to define  
554 a measurement session. This decision must be driven by the definition of short-term  
555 variation as a variation that is not “seen” by the selection process (see *e.g.* Price  
556 and Boag, 1987, p. 279 for a similar analogy), in contrast to persistent effects that  
557 are potentially under selection. This distinction ultimately depends on the trait, on  
558 the system under study and on the research question that is asked, because some  
559 traits may fluctuate on extremely short time scales (minutes or days), while others

560 remain constant across an entire adult's life.

561 The application to the snow vole data, where we varied the measurement session  
562 length from four days to one month, illustrated that longer measurement sessions  
563 automatically capture more variability, that is, the estimated error variance  $\hat{\sigma}_{e_m}^2$   
564 increased. Consequently, unreasonably long measurement sessions may lead to over-  
565 corrected estimates of the parameters of interest. On the other hand, considering  
566 measurement sessions that are too short may lead to an insufficient number of within-  
567 session repeats, or they may fail to identify transient variability that is biologically  
568 irrelevant. This makes clear that a careful definition of measurement session length  
569 is important already at the design stage of a study.

570 If one is uncertain whether repeated measurements capture effects relevant to se-  
571 lection or not, would averaging over repeats result in better estimates of quantitative  
572 genetic measures? Averaging methods have been proposed specifically to reduce bias  
573 that emerges due to measurement error and transient effects (Carbonaro et al., 2009;  
574 Zheng et al., 2016). While averaging will alleviate bias by reducing the error variance  
575 in the mean, it will not eliminate it completely. This can be seen from the fact that  
576 averaging over  $K$  within-session repeats for all animals and measurement sessions,  
577 the variance  $\sigma_{e_m}^2$  is reduced to  $\sigma_{e_m}^2 = \sigma_{e_m}^2 / K$ , assuming independence of the error  
578 term. Unless  $K$  is large,  $\sigma_{e_m}^2$  will not approach zero. Moreover, this practice only  
579 works if all animals have the same number of repeats within all measurement ses-  
580 sions, but it will not work in the unbalanced sampling design so common in studies  
581 of natural populations.

582 Our method approaches the problem of measurement error and transient fluc-  
583 tuations by assuming a dichotomous distinction between short-term and long-term  
584 repeats. An alternative perspective of within-animal repeated measurements could  
585 take a continuous view, recalling that repeated measurements are usually correlated,  
586 even when taken across long time spans, and that the correlation increases the closer  
587 in time the measurements were taken. A more sophisticated model could thus take  
588 into account that the residual component in the model changes continuously, and  
589 introduce a time-dependent correlation structure instead of simply distinguishing  
590 between short-term and long-term repeats. Such a model might be beneficial if  
591 repeats were not taken in clearly defined measurement sessions, although such a  
592 temporal correlation term introduces another level of model complexity, and thus  
593 entails other challenges.

594 It may sometimes not be possible to take multiple measurements on the same  
595 individual, or to repeat a measurement within a session. However, it may still be  
596 feasible to include an appropriate random effect in the absence of short-term repeats,  
597 provided that knowledge about the error variance is available, *e.g.* from previous  
598 studies that used the same measurement devices, from a subset of the data, or from

other “expert” knowledge. The Bayesian framework is ideal in this regard, because it is straightforward to include random effects with a very strong (or even fixed) prior on the respective variance component. Such Bayesian models provide error-aware estimates that are equivalent to those illustrated in Table 1, but with the additional advantage that posterior distributions naturally reflect all uncertainty that is present in the parameters, including the uncertainty that is incorporated in the prior distribution of the error variance.

Measurement error and transient fluctuations bias some, but not all quantitative genetic inferences. When  $\sigma_{em}^2 > 0$ , the naive estimates of  $h^2$ ,  $\beta_z$  and  $R_{BE}$  are attenuated by the same factor  $\lambda < 1$ , but other components, such as the selection differential  $S$  or  $R_{STS}$ , are not affected (Table 1). The robustness of the secondary theorem of selection to measurement error can certainly be seen as an advantage over the breeder’s equation. Nevertheless, the Robertson-Price identity does not model selection explicitly, and thus says little about the selective processes. The Robertson-Price equation can be used to check the consistency of predictions made from the breeder’s equation, but the breeder’s equation remains necessary to test hypothesis about the causal nature of selection (Morrissey et al., 2012; Bonnet et al., 2017). Another quantity that is unaffected by independent transient effects, which we however did not further elaborate on here, is *evolvability*, defined as the squared coefficient of variation  $I = \sigma_A^2 / \bar{z}^2$ , where  $\bar{z}$  denotes the mean phenotypic value (Houle, 1992). Evolvability is often used as an alternative to heritability, and is interpreted as the *opportunity for selection* (Crow, 1958). Not only  $\sigma_A^2$ , but also  $\bar{z}$  can be consistently estimated using  $z^*$ , namely because the expected values  $E[z^*] = E[z]$  due to the independence and zero mean of the error term. For completeness, we added evolvability to Table 1.

A critical assumption of our models was that the error components are independent of the phenotypic trait under study, but also independent of fitness or any covariates in the animal model or the selection model. While the small changes in  $\hat{R}_{STS}$  that we observed in the snow vole application with one-month measurement sessions could be due to pure estimation stochasticity, an alternative interpretation is that the measurement error in the data are not independent of the animal’s fitness. At least two processes could lead to a correlation between the measurement error in mass and fitness in snow voles. First, pregnant females will experience temporally increased body mass, and we expect the positive deviation from the true body mass to be correlated with fitness, because a pregnant animal is likely to have a higher expected number of offspring over its entire lifespan. And second, some of the snow voles were not fully grown when measured, and juveniles are more likely to survive if they keep growing, so that deviations from mean mass over the measurement session period would be non-randomly associated with life-time fitness.

So far, we have focused on traits that can change relatively quickly throughout the life of an individual, such as body mass, or physiological and behavioral traits. Traits that remain constant after a certain age facilitate the isolation of measurement error, because the residual variance term is then indistinguishable from the error term, given that a permanent environmental (*i.e.* individual-specific) effect is included in the model. In such a situation it is sufficient to estimate  $\sigma_R^2$ , which then automatically corresponds to the measurement error variance, while  $\sigma_{PE}^2$  captures all the environmental variability. However, not many traits will fit that description. The majority of traits, even seemingly stable traits such as skeletal traits, are in fact variable over time (Price and Grant, 1984; Smith et al., 1986).

We have shown that dealing appropriately with measurement error and transient fluctuations of phenotypic traits in quantitative genetic analyses requires the inclusion of additional variance components. Quantitative genetic analyses often differ in the variance components that are included to account for important dependencies in the data (Meffert et al., 2002; Palucci et al., 2007; Kruuk and Hadfield, 2007; Hadfield et al., 2013). Besides the importance of separating the right variance components, it has been widely discussed which of the components are to be included in the denominator of heritability estimates, although the focus has been mainly on the proper handling of variances that are captured by the fixed effects (Wilson, 2008; de Villemereuil et al., 2018). We hope that our treatment of measurement error in quantitative genetic analyses sparks new discussions of what should be included in the denominator when heritability is calculated.

The methods presented in this paper have been developed and implemented for continuous phenotypic traits. Binary, categorical or count traits may also suffer from measurement error, which is then denoted as misclassification error (Copas, 1988; Magder and Hughes, 1997; Küchenhoff et al., 2006), or as miscounting error (*e.g.* Muff et al., 2018). Models for non-Gaussian traits are usually formulated in a generalized linear model framework (Nakagawa and Schielzeth, 2010; de Villemereuil et al., 2016) and require the use of a link function (*e.g.* the logistic or log link). In these cases, it will often not be possible to obtain unbiased estimates of quantitative genetic parameters by adding an error term to the linear predictor as we have done here for continuous traits. Obtaining unbiased estimates of quantitative genetic parameters in the presence of misclassification and miscounting error will require extended modelling strategies, such as hierarchical models with an explicit level for the error process.

We hope that the concepts and methods provided here serve as a useful starting point when estimating quantitative genetics parameters in the presence of measurement error or transient, irrelevant fluctuations in phenotypic traits. The proposed approaches are relatively straightforward to implement, but further generalizations

677 are possible and will hopefully follow in the future.



## Supporting information:

**Appendix 1:** Supplementary text and figures (pdf)

**Appendix 2:** Supplementary text and figures for simulation study (pdf)

**Appendix 3:** R script for the simulation and analysis of pedigree data

**Appendix 4:** R script for heritability in snow voles

**Appendix 5:** R script for selection in snow voles

**Appendix 6:** R script for response to selection in snow voles.

## References

Bonnet, T., P. Wandeler, G. Camenisch, and E. Postma (2017). Bigger is fitter? Quantitative genetic decomposition of selection reveals an adaptive evolution decline of body mass in a wild rodent population. *PLOS Biology* 15, e1002592.

Carbonaro, F., T. Andrew, D. A. Mackey, T. L. Young, T. D. Spector, and C. J. Hammond (2009). Repeated measures of intraocular pressure result in higher heritability and greater power in genetic linkage studies. *Investigative Ophthalmology and Visual Science* 50, 5115–5119.

Carroll, R. J., D. Ruppert, L. A. Stefanski, and C. M. Crainiceanu (2006). *Measurement error in nonlinear models, a modern perspective*. Boca Raton: Chapman and Hall.

Charmantier, A., D. Garant, and L. E. B. Kruuk (2014). *Quantitative Genetics in the Wild*. Oxford: Oxford University Press.

Charmantier, A. and D. Reale (2005). How do misassigned paternities affect the estimation of heritability in the wild? *Molecular Ecology* 14, 2839–2850.

Copas, J. B. (1988). Binary regression models for contaminated data (with discussion). *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 50, 225–265.

Crow, J. F. (1958). Some possibilities for measuring selection intensities in man. *Human Biology* 30, 1–13.

de Boer, I. J. M. and I. Hoeschele (1993). Genetic evaluation methods for populations with dominance and inbreeding. *Theoretical Applied Genetics* 86, 245–258.

- 707 de Villemereuil, P., M. B. Morrissey, S. Nakagawa, and H. Schielzeth (2018). Fixed  
708 effect variance and the estimation of the heritability: Issues and solutions. *Journal*  
709 *of Evolutionary Biology* 31, 621–632.
- 710 de Villemereuil, P., H. Schielzeth, S. Nakagawa, and M. B. Morrissey (2016). General  
711 methods for evolutionary quantitative genetic inference from generalized mixed  
712 models. *Genetics* 204, 1281–1294.
- 713 Dohm, M. R. (2002). Repeatability estimates do not always set an upper limit to  
714 heritability. *Functional Ecology* 16, 273–280.
- 715 Falconer, D. S. and T. F. C. Mackay (1996). *Introduction to Quantitative Genetics*.  
716 Burnt Mill, Harlow, Essex, England: Pearson.
- 717 Fisher, R. A. (1930). *The Genetical Theory of Natural Selection*. Oxford, UK:  
718 Oxford University Press.
- 719 Fuller, W. A. (1987). *Measurement Error Models*. New York: John Wiley & Sons.
- 720 Ge, T., A. J. Holmes, R. L. Buckner, J. W. Smoller, and M. Sabuncu (2017). Her-  
721 itability analysis with repeat measurements and its application to resting-state  
722 functional connectivity. *PNAS* 114, 5521–5526.
- 723 Griffith, S. C., I. P. F. Owens, and K. A. Thuman (2002). Extrapair paternity  
724 in birds: a review of interspecific variation and adaptive function. *Molecular*  
725 *Ecology* 11, 2195–2212.
- 726 Hadfield, J. D. (2008). Estimating evolutionary parameters when viability selection  
727 is operating. *Proceedings of the Royal Society of London B: Biological Sciences*,  
728 *The Royal Society* 275, 723–734.
- 729 Hadfield, J. D., E. A. Heap, F. Bayer, E. A. Mittell, and N. M. Crouch (2013).  
730 Disentangling genetic and prenatal sources of familial resemblance across ontogeny  
731 in a wild passerine. *Evolution* 67, 2701–2713.
- 732 Henderson, C. R. (1976). A simple method for computing the inverse of a numerator  
733 relationship matrix used in prediction of breeding values. *Biometrics* 32, 69–83.
- 734 Hill, W. G. (2014). Applications of population genetics to animal breeding, from  
735 Wright, Fisher and Lush to genomic prediction. *Genetics* 196, 1–16.
- 736 Hoffmann, A. A. (2000). Laboratory and field heritabilities: Lessons from  
737 *Drosophila*. In T. Mousseau, S. B., and J. Endler (Eds.), *Adaptive Genetic Vari-*  
738 *ation in the Wild*. New York, Oxford: Oxford Univ Press.

- 739 Houle, D. (1992). Comparing evolvability and variability of quantitative traits.  
740 *Genetics* 130, 195–204.
- 741 Keller, L. F., P. R. Grant, B. R. Grant, and K. Petren (2001). Heritability of  
742 morphological traits in Darwin’s Finches: misidentified paternity and maternal  
743 effects. *Heredity* 87, 325–336.
- 744 Keller, L. F. and A. J. Van Noordwijk (1993). A method to isolate environmental  
745 effects on nestling growth, illustrated with examples from the Great Tit (*Parus*  
746 *major*). *Functional Ecology* 7, 493–502.
- 747 Kruuk, L. E. B. (2004). Estimating genetic parameters in natural populations using  
748 the ‘animal model’. *Philosophical Transactions of the Royal Society B: Biological*  
749 *Sciences* 359, 873–890.
- 750 Kruuk, L. E. B. and J. D. Hadfield (2007). How to separate genetic and environ-  
751 mental causes of similarity between relatives. *Journal of Evolutionary Biology* 20,  
752 1890–1903.
- 753 Küchenhoff, H., S. M. Mwalili, and E. Lesaffre (2006). A general method for dealing  
754 with misclassification in regression: The misclassification SIMEX. *Biometrics* 62,  
755 85–96.
- 756 Lande, R. and S. J. Arnold (1983). The measurement of selection on correlated  
757 characters. *Evolution* 37, 1210–1226.
- 758 Lush, J. L. (1937). *Animal breeding plans*. Ames, Iowa: Iowa State College Press.
- 759 Lynch, M. and B. Walsh (1998). *Genetics and Analysis of Quantitative Traits*.  
760 Sunderland, MA: Sinauer Associates.
- 761 Macgregor, S., B. K. Cornes, N. G. Martin, and P. M. Visscher (2006). Bias, precision  
762 and heritability of self-reported and clinically measured height in Australian twins.  
763 *Human Genetics* 120, 571–580.
- 764 Magder, L. S. and J. P. Hughes (1997). Logistic regression when the outcome is  
765 measured with uncertainty. *American Journal of Epidemiology* 146, 195–203.
- 766 Meffert, L. M., S. K. Hicks, and J. L. Regan (2002). Nonadditive genetic effects in  
767 animal behavior. *The American Naturalist* 160 Suppl 6, S198–S213.
- 768 Meuwissen, T. H. E., B. J. Hayes, and M. E. Goddard (2001). Prediction of total  
769 genetic value using genome-wide dense marker maps. *Genetics* 157, 1819–1829.

- 770 Mitchell-Olds, T. and R. G. Shaw (1987). Regression analysis of natural selection:  
771 statistical inference and biological interpretation. *Evolution* 41, 1149–1161.
- 772 Møller, A. and M. D. Jennions (2002). How much variance can be explained by  
773 ecologists and evolutionary biologists? *Oecologia* 132(4), 492–500.
- 774 Morrissey, M. B. and I. B. J. Goudie (2016). Analytical results for directional and  
775 quadratic selection gradients for log-linear models of fitness functions. *bioRxiv*.  
776 <https://www.biorxiv.org/content/early/2016/02/22/040618>.
- 777 Morrissey, M. B., L. E. B. Kruuk, and A. J. Wilson (2010). The danger of applying  
778 the breeder’s equation in observational studies of natural populations. *Journal of*  
779 *Evolutionary Biology* 23, 2277–2288.
- 780 Morrissey, M. B., D. J. Parker, P. Korsten, J. M. Pemberton, L. E. B. Kruuk, and  
781 A. J. Wilson (2012). The prediction of adaptive evolution: empirical application  
782 of the secondary theorem of selection and comparison to the breeder’s equation.  
783 *Evolution* 66, 2399–2410.
- 784 Morrissey, M. B. and K. Sakrejda (2013). Unification of regression-based methods  
785 for the analysis of natural selection. *Evolution* 67(7), 2094–2100.
- 786 Muff, S., M. A. Puhani, and L. Held (2018). Bias away from the Null due to mis-  
787 counted outcomes? A case study on the TORCH trial. *Statistical Methods in*  
788 *Medical Research*. In press.
- 789 Muff, S., A. Riebler, L. Held, H. Rue, and P. Saner (2015). Bayesian analysis  
790 of measurement error models using integrated nested Laplace approximations.  
791 *Journal of the Royal Statistical Society, Applied Statistics Series C* 64, 231–252.
- 792 Nakagawa, S. and H. Schielzeth (2010). Repeatability for Gaussian and non-  
793 Gaussian data: a practical guide for biologists. *Biological Reviews of the Cam-*  
794 *bridge Philosophical Society* 85, 935–956.
- 795 Palucci, V., L. R. Schaeffer, F. Miglior, and V. Osborne (2007). Non-additive ge-  
796 netic effects for fertility traits in Canadian Holstein cattle. *Genetics Selection*  
797 *Evolution* 39, 181–193.
- 798 Peek, M. S., A. J. Leffler, S. D. Flint, and R. J. Ryel (2003). How much variance is  
799 explained by ecologists? Additional perspectives. *Oecologia* 137(2), 161–170.
- 800 Price, G. R. (1970). Selection and covariance. *Nature* 227, 520–521.
- 801 Price, T. D. and P. T. Boag (1987). Selection in natural populations of birds. In  
802 F. Cooke, , and P. Buckley (Eds.), *Avian Genetics*, pp. 257 – 287. Academic Press.

- 803 Price, T. D. and P. R. Grant (1984). Life history traits and natural selection for  
804 small body size in a population of Darwin's Finches. *Evolution* 38, 483–494.
- 805 Richardson, S. and W. R. Gilks (1993). Conditional independence models for epi-  
806 demiological studies with covariate measurement error. *Statistics in Medicine* 12,  
807 1703–1722.
- 808 Robertson, A. (1966). A mathematical model of the culling process in dairy cattle.  
809 *Animal Science* 8, 95–108.
- 810 Roff, D. A. (2007). A centennial celebration for quantitative genetics. *Evolution* 61,  
811 1017–1032.
- 812 Rue, H., S. Martino, and N. Chopin (2009). Approximate Bayesian inference for  
813 latent Gaussian models by using integrated nested Laplace approximations (with  
814 discussion). *Journal of the Royal Statistical Society Series B (Statistical Method-*  
815 *ology)* 71, 319–392.
- 816 Senneke, S. L., M. D. MacNeil, and L. D. Van Vleck (2004). Effects of sire misiden-  
817 tification on estimates of genetic parameters for birth and weaning weights in  
818 Hereford cattle. *Journal of Animal Science* 82, 2307–2312.
- 819 Smith, J. N. M., P. Arcese, and D. Schuller (1986). Song sparrows grow and shrink  
820 with age. *AUK* 103, 210–212.
- 821 Steinsland, I., C. T. Larsen, A. Roulin, and H. Jensen (2014). Quantitative genetic  
822 modeling and inference in the presence of nonignorable missing data. *Evolution* 68,  
823 1735–1747.
- 824 Stephens, D. A. and P. Dellaportas (1992). Bayesian analysis of generalised linear  
825 models with covariate measurement error. In J. M. Bernardo, J. O. Berger, A. P.  
826 Dawid, and A. F. M. Smith (Eds.), *Bayesian Statistics 4*. Oxford Univ Press.
- 827 van der Sluis, S., M. Verhage, D. Posthuma, and C. V. Dolan (2010). Phenotypic  
828 complexity, measurement bias, and poor phenotypic resolution contribute to the  
829 missing heritability problem in genetic association studies. *PLOS One* 5, e13929.
- 830 Wilson, A. J. (2008). Why  $h^2$  does not always equal  $VA/VP$ ? *Journal of Evolutionary*  
831 *Biology* 21, 647–650.
- 832 Wilson, A. J., D. Réale, M. N. Clements, M. B. Morrissey, E. Postma, C. A. Walling,  
833 L. E. B. Kruuk, and D. H. Nussey (2010). An ecologist's guide to the animal model.  
834 *Journal of Animal Ecology* 79, 13–26.

- 835 Wolak, M. E. and J. M. Reid (2017). Accounting for genetic differences among  
836 unknown parents in microevolutionary studies: how to include genetic groups in  
837 quantitative genetic animal models. *Journal of Animal Ecology* 86, 7–20.
- 838 Zheng, Y., R. Plomin, and S. von Stumm (2016). Heritability of intraindividual  
839 mean and variability of positive and negative affect: genetic analysis of daily  
840 affect ratings over a month. *Psychological Science* 27, 1611–1619.

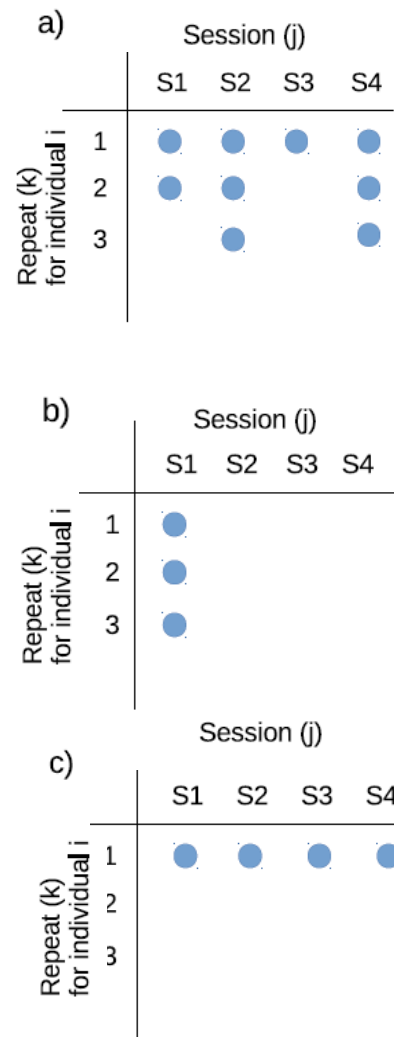


Figure 1: Schematic representation of three study designs, where one individual is measured a) multiple times across multiple measurement sessions, b) multiple times in one single measurement session, or c) one single time across multiple measurement sessions. Only case a) allows to disentangle the measurement error variance  $\sigma_{e_m}^2$  and the permanent environmental effects  $\sigma_{PE}^2$  from  $\sigma_R^2$ , while case b) allows to separate only the measurement error variance and case c) only allows to disentangle permanent environmental effects.



841

## Tables

842

Parameter	Effect of ME	Biased parameter
$\sigma_A^2$	unbiased	-
$\sigma_{PE}^2$	unbiased	-
$\sigma_R^2$	biased	$\sigma_R^2 + \sigma_e^2$
$h^2$	biased	$\lambda h^2$
$\beta_z$	biased	$\lambda \beta_z$
$\sigma_p(\mathbf{z}, \mathbf{w}) = S$	unbiased	-
$\sigma_a(\mathbf{z}, \mathbf{w}) = R_{STS}$	unbiased	-
$R_{BE}$	biased	$\lambda R_{BE}$
$I$	unbiased	-

Table 1: Overview of the effects of measurement error and transient fluctuations (ME) in a quantitative trait on important quantitative genetic parameters. The table indicates for each parameter whether it is biased or unbiased. For biased parameters the quantities are given that are estimated when ignoring transient effects in the quantitative genetic models.  $\lambda$  is the reliability ratio, defined as  $\lambda = \frac{\sigma_P^2}{\sigma_P^2 + \sigma_{e_m}^2}$ . For notation see the main text.

model	$\hat{h}^2$	$\hat{\sigma}_A^2$	$\hat{\sigma}_{PE}^2$	$\hat{\sigma}_M^2$	$\hat{\sigma}_R^2$	$\hat{\sigma}_{em}^2$
naive	0.14 [0.07, 0.25]	3.40 [1.41, 6.15]	6.09 [4.33, 8.51]	1.16 [0.56, 2.84]	12.40 [11.78, 13.21]	-
error-aware (4-day measurement session)	0.23 [0.09, 0.33]	3.97 [1.46, 6.06]	5.62 [3.68, 7.68]	1.48 [0.57, 2.73]	6.58 [5.76, 7.82]	6.07 [5.54, 7.05]
error-aware (one-month measurement session)	0.24 [0.10, 0.37]	3.82 [1.17, 5.84]	4.78 [3.16, 7.21]	1.58 [0.61, 2.86]	5.77 [4.78, 6.71]	7.91 [7.15, 8.38]

Table 2: Estimates of quantitative genetic parameters of body mass in snow voles using naive and error-aware models. The posterior modes of variance components and heritability are given, together with their 95% credible intervals (in brackets).

model	$\hat{\beta}_z$	$p$ -value
naïve	0.065	< 0.001
error-aware (4-day measurement session)	0.104	< 0.001
error-aware (one-month measurement session)	0.104	< 0.001

Table 3: Estimates of selection gradients ( $\hat{\beta}_z$ ) for body mass in snow voles, derived from naïve (ML estimate) and error-aware models (posterior means). For both types of models, Bayesian  $p$ -values were derived from zero-inflated Poisson regressions.

model	$\hat{R}_{\text{STS}}$	95% CI	$\hat{R}_{\text{BE}}$	95% CI
naive	-0.17	[-0.54, 0.18]	0.10	[0.05, 0.17]
error-aware (4-day measurement session)	-0.17	[-0.51, 0.19]	0.16	[0.06, 0.23]
error-aware (one-month measurement session)	-0.14	[-0.53, 0.17]	0.17	[0.07, 0.26]

Table 4: Response to selection for body mass in snow voles (posterior modes and 95% credible intervals) estimated with the breeder’s equation ( $\hat{R}_{\text{BE}}$ ) and with the secondary theorem of selection ( $\hat{R}_{\text{STS}}$ ). Results are shown for the naive and the error-aware models.